UNIBP: TOWARD UNIVERSAL BACKDOOR PURIFICATION VIA FINE-TUNING

Anonymous authors
Paper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026027028

029

031

033

034

037

038

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Deep neural networks (DNNs) remain vulnerable to backdoor attacks, perpetuating an arms race between attacks and defenses. Despite their efficacy against classical threats, mainstream defenses often fail under more advanced, defense-aware attacks, particularly clean-label variants that can evade decision-boundary shifting and neuron-pruning defenses. We present UniBP, a universal post-training defense that operates with only 1% of the original training data and unveils the relationship between batch normalization (BN) behavior and backdoor effects. At a high level, UniBP scrutinizes BN layers' affine parameters and statistics using a small clean subset (i.e., as small as 1% of the training data) to find the most impactful affine parameters for reactivating the backdoor, then prunes them and applies masked fine-tuning to remove the backdoor effects. We compare our method against 5 SOTA defenses, 5 backdoor attacks, and various attack/defense conditions, and show that UniBP consistently reduces the attack success rate from more than 90% to less than 5% while preserving clean performance, whereas other baselines degrade under smaller fine-tuning sets or stronger poisoning techniques. Our code is publicly available at https://anonymous.4open.science/r/ UniBP-BackdoorPostDefense/README.md.

1 Introduction

Deep neural networks (DNNs) have achieved remarkable success across a wide range of applications, including image classification, speech recognition, and natural language processing (Mienye & Swart, 2024; Samek et al., 2021; Noor & Ige, 2025). However, their vulnerability to backdoor attacks has raised serious concerns about their robustness in security-critical settings (Li et al., 2022; 2023b; Zhang et al., 2024; Wan et al., 2024; Cheng et al., 2025). In a backdoor attack, an adversary injects malicious patterns, which are commonly referred to as triggers into the training data. As a result, the model performs normally on clean inputs but misclassifies inputs containing the trigger in a controlled manner.

Backdoor attacks. Backdoor strategies have continued to evolve, becoming increasingly stealthy and effective. Early dirty-label methods such as BadNets (Gu et al., 2019) poison both inputs and labels, while later attacks like WaNet (Nguyen & Tran, 2021) apply subtle, visually faithful transformations that embed nearly-invisible triggers. More recent adaptive variants, including COMBAT (Huynh et al., 2024) and SBL (Sequential Learning Generates Resilient Backdoors) (Pham et al., 2024a), are explicitly crafted to bypass existing defenses, for example, by operating in clean-label regimes or by manipulating training dynamics to produce resilient, detection-aware backdoors. These advancements challenge traditional defense paradigms.

Defenses. In response, the literature spans adversarial training, input sanitization, and post-training defense. Recent methods are more focusing on the latest approach due to its practiacability in the erea of transfer learning, and where the training phase is not intervented (Min et al., 2024; Lin et al., 2024). epresentative methods include Neural Cleanse (Wang et al., 2019) and STRIP (Gao et al., 2019), which serve as post-training defenses: Neural Cleanse reverse-engineers class-wise minimal triggers to expose anomalies, and STRIP perturbs inputs and measures prediction entropy to detect triggered samples at inference. More recent defenses such as NAD (Li et al., 2021c), I-BAU (Zeng et al., 2021), ANP (Wu & Wang, 2021), and FST (Min et al., 2024) aim to handle a broader range of attacks using a clean dataset. They respectively distill clean behavior from a teacher (NAD), unlearn backdoors

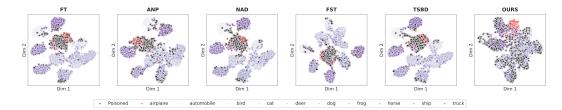


Figure 1: **t-SNE of feature embeddings on CIFAR-10.** Projection of penultimate-layer features for the backdoored model (PRETRAINED) and after applying defenses (ANP, FST, FT, NAD, TSBD, and UniBP). Clean samples are colored by class; poisoned samples are shown in black. All baselines fail in disrupting the overlapping representation of the backdoored data and the clean data of the targeted class (red).

with a minimax objective on a small clean set (I-BAU), prune adversarially sensitive neurons (ANP), and fine-tune to separate backdoor from clean features (FST). However, these defenses still primarily target dirty-label settings, and empirically we find that they are ineffective or unstable against newer attacks such as COMBAT (cf. Figure 1), especially under varying attacker/defender configurations and data budgets.

Our approach. In this paper, we present a universal and practical *post-training* defense grounded in a key observation: Batch Normalization (BN) layers encode distributional statistics of both clean and poisoned data, and backdoor behavior exploits these statistics to steer specific activation pathways. Our method (i) *rectifies and aligns* the backdoored model's BN statistics during fine-tuning to find channels most responsible for trigger activation, then (ii) *resets* a targeted subset of BN affine parameters and (iii) applies *masked-gradient fine-tuning* to prevent reactivation by a malicious trigger. This yields effective purification of pretrained models *without* prior knowledge of attack type, trigger pattern, or poisoned locations, and operates with minimal clean data and assumptions. In practice, the procedure is effective across various backdoor attacks, stable across diverse attack conditions, and architecture-agnostic.

To summarize, our main contributions are as follows:

- We unveil the relationship of BN layers' affine parameters and statistics toward the backdoor effect, and show that only a subset (i.e., 0.01%) of these parameters can sustainably disrupt the backdoor's attack success rate.
- We then introduce UniBP, a post-training defense that finds these affine parameters, then conducts pruning and masked fine-tuning to remove the backdoor from a poisoned model.
- We empirically show that prior fine-tuning defenses are often ineffective and unstable across major backdoor families. In contrast, our method is universal in that it is consistently effective against traditional (BadNets (Gu et al., 2019), WaNet (Nguyen & Tran, 2021)), clean-label (LC (Turner et al., 2019)), and adaptive (COMBAT (Huynh et al., 2024), SBL (Pham et al., 2024b)) backdoor attacks.
- We rigorously evaluate UniBP across a swath of attack settings and model architectures. We show that UniBP (1) preserves clean accuracy while maintaining stability and resilience against each attack, (2) requires only a small amount of clean data, and (3) requires *no* assumptions about the implanted backdoor.

2 Related Works

2.1 Backdoor Attacks

Backdoor attacks aim to mislead a victim model into predicting the target label when a trigger is present in the input while keeping the model performance unchanged on clean data. Backdoor attacks are categorized into *dirty-label* (Chen et al., 2017; Li et al., 2021b; Wang et al., 2022) and *clean-label* (Barni et al., 2019; Ning et al., 2021; Zeng et al., 2023) based on whether the attacker crafts a trigger in a way that changes the underlying label of the poisoned image. In *dirty-label backdoor attacks*, while the seminal work, Badnet (Gu et al., 2019), uses a single or a pattern of bright pixels as a trigger, later works have focused on making the trigger undetectable, e.g. by using image

warping (Nguyen & Tran, 2021). The inconsistency between image and label in dirty-label backdoor attacks is often visually detectable by humans. In *clean-label backdoor attacks* (Barni et al., 2019; Ning et al., 2021; Zeng et al., 2023), triggers are only added to the data in the target class. LC (Turner et al., 2019) perturbs the original input to ensure that the model learns the trigger. COMBAT (Huynh et al., 2024) aims to learn an effective trigger generator by adopting an alternate training process that optimizes the generator and a surrogate model. Building on both types of attacks, recent work seeks to improve resilience of the backdoored model against fine-tuning, e.g., SBL (Pham et al., 2024b) traps the backdoored model within the backdoored region via continual learning. Specifically designed to resist fine-tuning defenses, SBL simulates the effects of fine-tuning during its second stage, allowing the backdoor to remain highly effective even after mainstream defenses are applied.

2.2 BACKDOOR DEFENSES

In response to the growing threat of backdoor attacks, various defensive techniques have been proposed that operate during two stages of model training: (1) training-stage and (2) post-training defenses. Training-stage defenses (Huang et al., 2022) aim to train a clean model even when the training data has been poisoned by an attacker. ABL (Li et al., 2021b) first isolates the backdoored data and then unlearns the isolated data using gradient ascent. D-ST/D-BR (Chen et al., 2022) leverages the insight that poisoned data are more sensitive to transformation compared to clean data, so they train a secure model from scratch or unlearn poisoned samples in a backdoored model. In contrast, Post-training defenses (Zheng et al., 2022a; Chen et al., 2018; Nguyen et al., 2024) aim to mitigate the backdoor effect on a poisoned model using a small set of known-clean data. This is usually achieved through pruning or fine-tuning. ANP (Wu & Wang, 2021) prunes sensitive neurons under adversarial neuron perturbation as they are likely to be related to the injected backdoor. NAD (Li et al., 2021c) introduces an attention distillation method which uses a teacher network to guide the fine-tuning of the backdoored network. FST (Min et al., 2024) encourages discrepancy between fine-tuned model and the original model to achieve feature shifts. Recently, TSBD (Lin et al., 2024) leverages the insight that neuron weight changes are highly-correlated in poisoned unlearning and clean unlearning, and (1) reinitialize neurons based on weight changes, and (2) fine-tune the model based on neuron activeness. PBP (Nguyen et al., 2024) first generates a neuron mask, then uses masked gradient optimization to eliminate backdoor effects. However, current SOTA defenses have not effectively tackled the newly proposed resilient backdoor attacks, including SBL and COMBAT, underscoring the need for more robust defense mechanisms.

3 METHODOLOGY

3.1 PROBLEM STATEMENTS

Much of the backdoor attack literature (Gu et al., 2017; Zheng et al., 2022b;a; Wang et al., 2023) assumes an "Outsourced Training Attack," where adversaries control training and users rely only on a held-out validation set. However, since backdoored models maintain high clean-data performance, validation alone is insufficient to verify whether a model is backdoored. To address this challenge, we adopt a defense setting where the defender acquires a backdoored model from an untrusted source and assumes access to a small subset of clean training data for fine-tuning \mathcal{D}_{ft} (Li et al., 2023a; 2021c). Backdoor defense/purification aims to eliminate the backdoor trigger while maintaining the model's performance on clean samples. This approach is particularly relevant when training data is no longer fully accessible due to retention or privacy policies.

Attacker's goals. Similar to most backdoor poisoning settings, we assume the attacker's goal is to alter the training procedure by using a small poisoned set, such that the resulting trained backdoored classifier, f_{θ^*} , differs from a cleanly trained classifier. An ideal f_{θ^*} has the same response to clean samples, whereas it generates an adversarially chosen prediction, $\tau(y)$, when applied to backdoored inputs, $\varphi(x)$.

Defender's goal. In contrast to the attacker, the defender—who has full access to the poisoned model f_{θ^*} and a limited benign fine-tuning set \mathcal{D}_{ft} to get a clean/purified model $f_{\hat{\theta}}$ must (1) remove backdoors from f_{θ^*} to ensure correct behavior on triggered inputs and (2) preserve the model's performance on normal inputs during purification. In this work, following related post-training defenses Min et al. (2024); Wang et al. (2023); Lin et al. (2024), we adopt the following assumptions in a compact form: (i) the defender has no information about the backdoor trigger or the adversary's

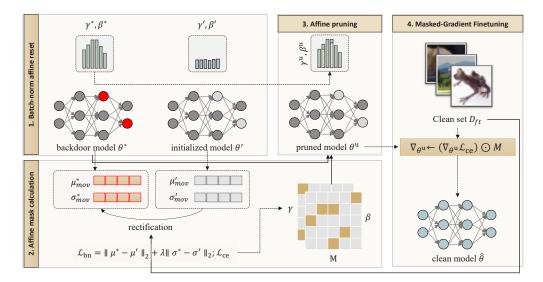


Figure 2: **UniBP** includes four phases. (1) *Batch-norm affine reset* reinitializes γ , β of the backdoored model θ^* to obtain θ' . (2) *Affine-mask calculation* rectifies BN moving statistics (μ_{mov} , σ_{mov}) and learns a selection mask M via the BN rectification loss $\mathcal{L}_{\text{bn}} = \|\mu' - \mu\|_2 + \lambda \|\sigma' - \sigma\|_2$ (with \mathcal{L}_{ce}). (3) *Affine pruning* removes suspect channels/affines, yielding θ^u with γ^u , β^u . (4) *Masked-gradient finetuning* on a small clean set D_{ft} updates only unmasked parameters $(\nabla_{\theta^u}\mathcal{L}_{\text{ce}}) \odot M$, producing the purified model $\hat{\theta}$.

accessibility (e.g., poisoning rate, insertion mechanism), and we make no assumptions about any trigger/watermark; (ii) the defender has no access to the original training procedure and cannot obtain the full training dataset to retrain a new model; and (iii) the defender can collect or access a small, clean dataset representative of the training distribution (covering all classes), and may combine it with any available portion of the training data. This setting aligns with common post-training defenses (Min et al., 2024; Wang et al., 2023).

3.2 RELATIONSHIP OF BN LAYERS AND BACKDOOR EFFECT.

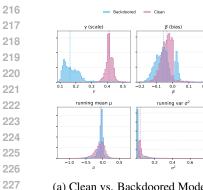
Finding 3.1: Backdoors shift BatchNorm (BN) statistics and affine parameter distributions

Training with a backdoor induces consistent, layer-dependent shifts in BN running means/variances and alters the distribution of BN affine parameters (γ, β) relative to clean baselines.

BatchNorm layers are often used in deep neural networks for the purposes of stabilizing and accelerating training (by reducing internal covariate shift), permitting larger learning rates, improving generalization via a mild regularization effect, and offering per-channel control through learnable affine parameters. Given a mini-batch of feature maps $x_{n,c,h,w}$ with batch size N and spatial size $H \times W$, BN computes:

$$\mu_{c} = \frac{1}{NHW} \sum_{n,h,w} x_{n,c,h,w}, \quad \sigma_{c}^{2} = \frac{1}{NHW} \sum_{n,h,w} \left(x_{n,c,h,w} - \mu_{c} \right)^{2}, \quad \hat{x}_{n,c,h,w} = \frac{x_{n,c,h,w} - \mu_{c}}{\sqrt{\sigma_{c}^{2} + \varepsilon}}$$
(1)

and outputs the affine-transformed activations as $y_{n,c,h,w} = \gamma_c \hat{x}_{n,c,h,w} + \beta_c$, where γ_c and β_c are learned affine (scale/shift) parameters for channel c, and ε ensures numerical stability. During training, (μ_c, σ_c^2) are computed from the current mini-batch while exponential moving averages are accumulated; at inference, these running estimates replace batch statistics. Our key insight (see 3.1) is that BN layers encode the training distribution via their running moments and affine parameters ??, and inserting a backdoor unavoidably shifts the distribution of the BN layers' statistic and affine parameters (see Figure 3a). Building on this observation, we articulate our second finding (3.2), which is central to our methodology: backdoor activation is governed by a small subset of BN affine channels; consequently, identifying and selectively editing these channels serves as an surprisingly effective lever for backdoor mitigation (cf. Figure 3c).



229

230

231

232

233 234 235

236 237

238

239

240 241

242 243

244

245

246

247

249

250

251

252

253

254

259

260

261

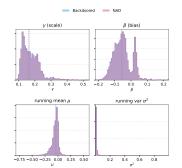
262

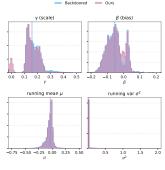
264

265 266

267

268





- (a) Clean vs. Backdoored Models
- (b) Backdoored vs. NAD Models
- (c) Backdoored vs. Our Models

Figure 3: BatchNorm statistics (μ, σ) and affine parameters (γ, β) for four model variants—clean, backdoored, NAD (Li et al., 2021a), and Ours—illustrating how backdoor training and purification affect BN layers. NAD leaves the backdoored BN statistics largely unchanged, whereas our method slightly shifts them while successfully removing the backdoor. ASR: clean 0.67%, backdoored 80.66%, NAD 78.66%, Ours 7.04% (lower is better).

Finding 3.2: Backdoor activation is bottlenecked by a sparse subset of BN affine parameters

Claim. A small fraction of BN affine channels (γ, β) disproportionately governs trigger activation; selectively perturbing or resetting these top-ranked channels sharply reduces ASR with minimal impact on clean accuracy.

UNIBP: DETAILED DESCRIPTION 3.3

High-Level Idea. Motivated by the two findings mentioned above, we introduce a defense method including four components. (i) batch-norm affine reset to create an initialized model θ' from the backdoored model θ^* ; (ii) affine mask calculation by calculating FIM while the initialized model is trained with rectification to align the BN stats with the backdoored model; (iii) this mask will be used to prune the corresponding highly influential neurons to remove the backdoor effect, achieving a pruned model θ^u ; (iv) this pruned model is then fine-tuned using masked-gradient training with a clean dataset to achieve the purified version $\hat{\theta}$.

Batch-norm affine reset. Given a backdoored model θ^* , we obtain the corresponding *re-initialized* model θ' by resetting BatchNorm affine parameters. Let \mathcal{B} be the set of BN layers in θ^* , and for each $\ell \in \mathcal{B}$ with C_{ℓ} channels let $(\gamma_{\ell}, \beta_{\ell}) \in \mathbb{R}^{C_{\ell}} \times \mathbb{R}^{C_{\ell}}$ denote its affine parameters (if present). For fixed reinit constants (γ_0, β_0) (i.e., which are set default as (1,0)), we define the operator \mathcal{R}_{BN} :

$$\theta' = \mathcal{R}_{\mathrm{BN}}(\theta^*; \gamma_0, \beta_0), \quad (\gamma'_\ell, \beta'_\ell) = \begin{cases} (\gamma_0 \, \mathbf{1}_{C_\ell}, \, \beta_0 \, \mathbf{1}_{C_\ell}) & \text{if the BN layer ℓ has affine parameters,} \\ (\gamma_\ell, \beta_\ell) & \text{otherwise.} \end{cases}$$

Affine Mask Calculation. From the initialized model θ' , we compute an *importance score* for each BN affine parameter that quantifies its contribution to rectifying the BatchNorm statistics of $\theta'(\mu'_{\ell}, \nu'_{\ell})$ toward those of the backdoored model $(\mu_{\ell}^*, v_{\ell}^*)$. This procedure mimics the alignment in which the statistics induced by a small clean fine-tuning set $\mathcal{D}_{\mathrm{ft}}$ are drawn toward the mixed (clean and poisoned) distribution used to train θ^* . To achieve this goal, we fine-tune the reinitialized model θ' by minimizing the rectification objective, and we quantify per-parameter importance via the (empirical) Fisher information computed on $\mathcal{D}_{\mathrm{ft}}$. Specifically, we use $\mathcal{L}_{\mathrm{rectify}}$ for optimization and estimate the diagonal Fisher for each parameter ϕ_i as in equation 3.

Let \mathcal{B} be the set of BN layers, for each $\ell \in \mathcal{S}$, let (μ'_{ℓ}, u'_{ℓ}) denote the per-channel batch mean/variance computed on the current mini-batch as in Equation 1, and let $(\mu_{\ell}^*, v_{\ell}^*)$ be the corresponding references from the backdoored model. We define the per-layer deviation loss function as follows:

$$\mathcal{L}_{\mathrm{BN}}^{(\ell)} = \left\| \hat{\mu}_{\ell} - \mu_{\ell}^* \right\|_2 + \lambda \left\| \hat{v}_{\ell} - v_{\ell}^* \right\|_2, \qquad \lambda = 0.05.$$

Then, the BN regularizer is calculated as: $\mathcal{L}_{\mathrm{BN}} = \frac{1}{|\mathcal{S}|} \sum_{\ell \in \mathcal{S}} \mathcal{L}_{\mathrm{BN}}^{(\ell)}$. This regularizer encourages the network's intermediate distributions to align with the reference (backdoored) normalization statistics, stabilizing activations without directly constraining (γ, β) . We then define the rectification objective by:

$$\mathcal{L}_{rectify} := \mathcal{L}_{CE}(x, y) + \log \mathcal{L}_{BN}. \tag{2}$$

Let Θ denote all trainable parameters and $\Theta_{\rm BN} \subset \Theta$ the set of BN affine entries $\{\gamma_{\ell,c},\beta_{\ell,c}:\ell\in\mathcal{B},\ 1\leq c\leq C_\ell\}$. We quantify per-parameter sensitivity under the rectification objective $\mathcal{L}_{\rm rectify}$ via the empirical (diagonal) Fisher:

$$\widehat{F}_{\theta_i}^{(\text{rect})} = \frac{1}{|\mathcal{D}_{\text{ft}}|} \sum_{(x,y) \in \mathcal{D}_{\text{ft}}} \|\nabla_{\theta_i} \mathcal{L}_{\text{rectify}}(x,y)\|^2, \qquad \theta_i \in \Theta.$$
 (3)

For BN affines we set the importance score $s_j := \widehat{F}_{\theta_j}^{(\mathrm{rect})}$ for each $\theta_j \in \Theta_{\mathrm{BN}}$.

Mask Construction. Let $K \in \mathbb{N}$ be the pruning budget (optionally $K = \lfloor r | \Theta_{\mathrm{BN}} | \rfloor$ for a ratio $r \in (0,1)$), and let τ be the K-th largest value of $\{s_j : \theta_j \in \Theta_{\mathrm{BN}}\}$. Define the binary mask $M_i \in \{0,1\}$ by

$$M_j = \mathbf{1}\{s_j < \tau\} = \begin{cases} 0, & \text{if } s_j \text{ is among the top-} K \text{ in } \Theta_{\text{BN}}, \\ 1, & \text{otherwise.} \end{cases}$$
 (4)

Affine Pruning. Pruning is one of the most popular methods to remove the effect of a subset of neurons on the model activation and prediction (Li et al., 2021a;c). To remove the backdoor effect, we prune the BatchNorm affine parameters whose corresponding mask values are zero. Concretely, for the k-th neuron , we set its weight $w_k=0$ if $M_k=0$ and keep it unchanged if $M_k=1$. Due to the binary masks, pruning is a discrete optimization problem that is difficult to solve within feasible time. To address this, we add a small Gaussian noise to the parameters at the pruned coordinates during fine-tuning. Given the BN affine parameters $\Theta=\{\theta_j\}$ and the affine mask M determined in the previous step, and $\Xi\sim\mathcal{N}(0,\sigma^2I)$ be i.i.d. noise. We use the masked-and-noised parameters:

$$\theta^{u} := \tilde{\Theta} = M \odot \Theta + (1 - M) \odot \Xi. \tag{5}$$

Masked-Gradient Finetuning. During this process, we zero out the gradient at the affine parameters which are pruned in the previous step. The objective for fine-tuning can be stated as follows:

$$\hat{\theta} := \min_{\theta} \mathbb{E}_{(\boldsymbol{x},y) \in \mathcal{D}_{\mathrm{ft}}} \mathcal{L}_{\mathrm{CE}}(f(\boldsymbol{x}; M \odot \theta^{u}), y), \qquad \nabla_{\theta} \leftarrow (\nabla_{\theta} \mathcal{L}_{\mathrm{CE}}) \odot M, \tag{6}$$

where $\hat{\theta}$ denotes the current parameters. The mask zeroes gradients only on BN-affine coordinates and leaves all other parameters trainable, preventing drift back toward the backdoored BatchNorm statistics while preserving clean behavior.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETTINGS

Attack Setup. We consider five distinct backdoor attack strategies: (1) BadNet (Gu et al., 2019), (2) Label Consistent (LC) (Turner et al., 2019), (3) Wanet (Nguyen & Tran, 2021), (4) COMBAT (Huynh et al., 2024), and (5) SBL (Pham et al., 2024a). BadNet and LC are two mainstream methods for dirty- and clean-label backdoor attacks, respectively, while COMBAT is the most recent optimized trigger clean-label backdoor attack. SBL is designed to preserve backdoor durability even when defenses are applied during fine-tuning. We leverage the BackdoorBench (Wu et al., 2022) framework using the authors' provided code for COMBAT and SBL to control trigger pattern, trigger size, and target label. We vary the poisoning rate from 1% to 10%. Unless otherwise stated, we adopt PreAct-ResNet-18 (He et al., 2016) and a 10% fine-tuning ratio by default. We evaluate on two benchmark datasets: CIFAR-10 (Krizhevsky et al., 2009) and GTSRB (Stallkamp et al., 2011). Owing to space constraints, we report a representative subset of results here; additional details and full results are provided in the Appendix.

Table 1: Comparison of SOTA defenses against multiple backdoor attacks with different fine-tuning ratios on the CIFAR-10 dataset with PreAct-ResNet18.

| Methods | Metrics | | | 0.1 | | | 0.05 | | | | | | | |
|------------|---------|--------|--------|--------|-------|-------|-------|--------|--------|--------|-------|-------|-------|--|
| Methods | | BadNet | LC | COMBAT | SBL | Wanet | Avg. | BadNet | LC | COMBAT | SBL | Wanet | Avg. | |
| Pretrained | C-ACC | 91.44 | 84.19 | 93.94 | 90.52 | 92.67 | 90.55 | 91.36 | 84.51 | 94.13 | 89.76 | 92.90 | 90.53 | |
| | ASR | 94.41 | 100.00 | 94.47 | 88.84 | 99.54 | 95.45 | 95.45 | 100.00 | 94.80 | 87.55 | 99.54 | 95.47 | |
| | DER | - | - | - | - | - | - | - | - | - | - | - | - | |
| | C-ACC | 90.56 | 90.00 | 93.46 | 90.79 | 92.50 | 91.46 | 88.69 | 90.51 | 94.01 | 89.46 | 92.33 | 91.00 | |
| FT | ASR | 1.47 | 17.53 | 72.83 | 83.85 | 13.91 | 37.92 | 2.22 | 100.00 | 96.17 | 89.92 | 14.97 | 60.66 | |
| | DER | 96.03 | 91.24 | 60.58 | 52.50 | 92.73 | 78.62 | 95.28 | 50.00 | 49.94 | 49.85 | 92.00 | 57.51 | |
| ANP | C-ACC | 83.51 | 79.17 | 85.18 | 88.77 | 83.62 | 84.05 | 84.40 | 84.51 | 92.14 | 84.48 | 84.83 | 86.07 | |
| | ASR | 0.00 | 6.65 | 7.58 | 0.04 | 0.02 | 2.86 | 0.02 | 100.00 | 88.81 | 62.48 | 0.00 | 50.26 | |
| | DER | 93.24 | 94.17 | 89.07 | 93.53 | 95.24 | 93.05 | 94.24 | 50.00 | 52.00 | 59.90 | 95.74 | 60.48 | |
| | C-ACC | 89.33 | 88.97 | 93.48 | 90.39 | 91.88 | 90.81 | 88.13 | 89.30 | 94.21 | 88.94 | 92.09 | 90.53 | |
| NAD | ASR | 2.08 | 18.43 | 70.96 | 64.80 | 9.98 | 33.25 | 2.81 | 59.06 | 97.66 | 73.08 | 1.83 | 46.89 | |
| | DER | 95.11 | 90.79 | 61.52 | 61.96 | 94.39 | 80.75 | 94.71 | 70.47 | 50.00 | 56.83 | 98.45 | 74.09 | |
| | C-ACC | 87.06 | 88.89 | 91.25 | 91.17 | 92.40 | 90.15 | 88.58 | 90.90 | 94.20 | 89.95 | 92.43 | 91.21 | |
| FST | ASR | 2.08 | 2.34 | 30.65 | 0.24 | 0.58 | 7.18 | 1.13 | 0.00 | 90.02 | 30.02 | 0.32 | 24.30 | |
| | DER | 93.98 | 98.83 | 80.57 | 94.30 | 99.35 | 93.41 | 95.77 | 1.00 | 52.39 | 78.77 | 99.38 | 65.46 | |
| | C-ACC | 90.13 | 89.06 | 92.91 | 91.43 | 92.48 | 91.20 | 90.00 | 90.74 | 92.28 | 88.14 | 92.43 | 90.72 | |
| TSBD | ASR | 1.78 | 15.16 | 35.57 | 84.68 | 1.08 | 27.65 | 2.12 | 93.00 | 81.64 | 79.20 | 1.29 | 51.45 | |
| | DER | 95.66 | 92.42 | 78.94 | 52.08 | 99.14 | 83.65 | 95.99 | 53.50 | 55.66 | 53.37 | 98.89 | 71.48 | |
| | C-ACC | 90.67 | 91.40 | 91.04 | 88.91 | 90.22 | 90.45 | 90.32 | 88.94 | 85.49 | 86.17 | 89.45 | 88.07 | |
| Ours | ASR | 1.12 | 2.50 | 10.28 | 2.18 | 4.74 | 4.16 | 4.91 | 5.08 | 7.30 | 4.84 | 2.64 | 4.95 | |
| | DER | 98.99 | 98.75 | 93.02 | 97.88 | 96.30 | 96.99 | 96.86 | 97.46 | 91.56 | 95.24 | 96.82 | 95.59 | |

Baselines. We consider five state-of-the-art defenses, representing a range of strategies aimed at mitigating the impact of backdoor attacks, from continued training on clean data to model pruning and reinitialization. These defenses include Fine-tuning (FT), NAD (Li et al., 2021a), ANP (Wu & Wang, 2021), FST (Min et al., 2024), and TSBD (Lin et al., 2024). We follow the suggested parameters from BackdoorBench.

Metrics. Following (Lin et al., 2024; Min et al., 2023; Zhu et al., 2023), we report *C-ACC* (clean accuracy), *ASR* (attack success rate), and $DER \in [0,1]$, which balances ASR reduction against utility: $DER = \frac{\max(0, \Delta ASR) - \max(0, \Delta C - ACC) + 1}{2}$, where ΔASR and ΔACC are the drop in ASR and C-ACC after applying defense on the backdoored model, respectively. We expect a good defense to have a large C-ACC, DER, and a small ASR. Following (Zeng et al., 2022), we mark [ASR] when ASR > 10% and [C-ACC] when C-ACC decreases by > 10%. We highlight the best among the six baselines with [DER]. We tag UniBP with [DER] when it is comparable to or better than the best baseline, where "comparable" means C-ACC gap < 2% and ASR gap < 4%.

4.2 MAIN RESULTS

We compare the performance of our method to five other defenses against five representative backdoor attacks. In this section, we present the main results on CIFAR-10 and GTSRB with a 10% poisoning ratio on PreAct-ResNet18 for illustration, which is shown in Table 1 and Table 2.

Performance of backdoor defenses on CIFAR-10 dataset. In the CIFAR-10 dataset, our method demonstrates consistent performance across different fine-tuning ratios and outperforms all state-of-the-art defenses on average. At a fine-tuning ratio of 0.1, where the defender can access relatively more clean data, ANP, FST, and our approach are all able to reduce the attack success rate (ASR) while maintaining high clean accuracy (C-ACC). Among them, our method achieves the highest average DER of 96.30%. In contrast, NAD and TSBD already show clear deficiencies against stronger attacks such as COMBAT and SBL, which are either input-dependent or explicitly designed to resist fine-tuning. When the fine-tuning ratio is reduced to 0.05, ANP fails to mitigate several attacks, with ASRs of 100.00% in LC, 88.81% in COMBAT and 62.48% in SBL, while FST achieves an ASR of 90.02% against COMBAT and 30.02% against SBL. By comparison, our method continues to maintain the most effective defense against all attacks, achieving an average C-ACC of 90.45% and an average ASR of only less than 5%.

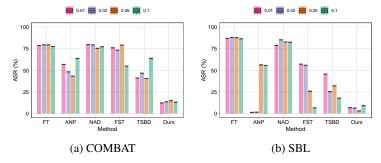


Figure 4: Defense results (ASR) under various fine-tuning ratio settings with COMBAT and SBL attacks. The experiments are conducted on the CIFAR-10 dataset.

Performance of backdoor defenses on GTSRB dataset. On the GTSRB dataset, our method follows a similar trend to that observed on CIFAR-10, consistently outperforming all SOTA defenses and achieving the highest average DER. At a fine-tuning ratio of 0.1, NAD and TSBD again show deficiencies, while ANP and FST also fail to effectively reduce ASR under stronger attacks. ANP achieves a 69.37% ASR and FST achieves 21.65% ASR against COMBAT, illustrating that these approaches struggle when the dataset becomes more complex. The same pattern is evident at a lower fine-tuning ratio of 0.05, where most defenses fail to mitigate at least one attack. *In contrast, our method maintains robust performance across all scenarios, achieving the highest DER of 95.43% for both fine-tuning scenarios.*

Table 2: Comparison of SOTA defenses against multiple backdoor attacks with different fine-tuning ratios on the GTSRB dataset with PreAct-ResNet18.

| Methods | Metrics | | | 0.1 | | | 0.05 | | | | | | | |
|------------|---------|--------|-------|--------|-------|-------|-------|--------|-------|--------|-------|-------|-------|--|
| memous | | BadNet | LC | COMBAT | SBL | Wanet | Avg. | BadNet | LC | COMBAT | SBL | Wanet | Avg. | |
| Pretrained | C-ACC | 96.85 | 92.45 | 99.07 | 97.36 | 96.19 | 96.38 | 96.94 | 92.41 | 97.97 | 97.29 | 97.45 | 96.41 | |
| | ASR | 94.29 | 99.24 | 69.53 | 91.96 | 99.53 | 90.91 | 94.61 | 99.96 | 76.07 | 90.97 | 99.14 | 92.15 | |
| | DER | - | - | - | - | - | - | - | - | - | - | - | - | |
| | C-ACC | 97.58 | 97.27 | 98.97 | 97.37 | 98.65 | 97.97 | 97.78 | 97.36 | 98.57 | 97.41 | 98.93 | 98.01 | |
| FT | ASR | 67.35 | 97.82 | 65.69 | 90.18 | 34.99 | 71.21 | 57.86 | 99.06 | 78.18 | 89.27 | 80.17 | 80.91 | |
| | DER | 63.47 | 50.71 | 51.87 | 50.89 | 82.27 | 59.84 | 68.38 | 50.45 | 50.00 | 50.85 | 59.49 | 55.83 | |
| ANP | C-ACC | 95.97 | 91.55 | 98.73 | 95.14 | 98.33 | 95.94 | 93.66 | 92.17 | 97.87 | 95.16 | 93.60 | 94.49 | |
| | ASR | 13.64 | 1.31 | 69.37 | 1.07 | 0.00 | 17.08 | 0.00 | 35.34 | 54.02 | 0.00 | 0.00 | 17.87 | |
| | DER | 89.89 | 98.52 | 49.91 | 94.34 | 99.76 | 86.48 | 95.67 | 82.19 | 60.98 | 94.42 | 97.65 | 86.18 | |
| NAD | C-ACC | 97.61 | 97.37 | 99.05 | 97.43 | 98.77 | 98.05 | 97.86 | 96.16 | 98.49 | 97.67 | 99.04 | 97.84 | |
| | ASR | 25.53 | 0.37 | 65.23 | 62.59 | 23.60 | 35.46 | 8.44 | 0.40 | 75.15 | 87.72 | 64.91 | 47.32 | |
| | DER | 84.38 | 99.44 | 52.14 | 64.69 | 87.97 | 77.72 | 93.09 | 99.78 | 50.46 | 51.63 | 67.12 | 72.41 | |
| | C-ACC | 94.50 | 97.28 | 97.71 | 97.22 | 98.89 | 97.12 | 93.64 | 95.17 | 97.81 | 97.33 | 98.66 | 96.52 | |
| FST | ASR | 0.00 | 0.00 | 21.65 | 0.00 | 0.00 | 4.33 | 0.00 | 0.00 | 63.15 | 0.00 | 0.00 | 12.63 | |
| | DER | 95.97 | 99.62 | 73.35 | 95.91 | 99.76 | 92.92 | 95.66 | 99.98 | 56.38 | 95.49 | 99.57 | 89.41 | |
| | C-ACC | 98.20 | 97.34 | 99.21 | 97.29 | 94.45 | 97.30 | 97.98 | 96.39 | 98.39 | 96.65 | 86.66 | 95.21 | |
| TSBD | ASR | 0.00 | 0.16 | 66.58 | 45.42 | 0.21 | 22.47 | 0.03 | 0.80 | 53.99 | 13.22 | 0.00 | 13.61 | |
| | DER | 97.15 | 99.29 | 51.48 | 73.24 | 98.79 | 83.99 | 97.29 | 99.58 | 61.04 | 88.56 | 94.18 | 88.13 | |
| | C-ACC | 97.43 | 98.22 | 97.38 | 97.10 | 95.19 | 97.06 | 97.86 | 97.40 | 90.63 | 96.27 | 97.16 | 95.86 | |
| Ours | ASR | 0.00 | 0.04 | 1.23 | 0.08 | 0.02 | 0.27 | 0.01 | 0.36 | 4.30 | 0.08 | 0.00 | 0.95 | |
| | DER | 97.15 | 99.60 | 83.31 | 95.81 | 99.26 | 95.43 | 97.30 | 99.87 | 74.77 | 94.99 | 99.43 | 93.27 | |

4.3 ABLATION STUDIES

In this section, we study the performance of different defenses under varied adversary ability and defender capability. Specifically, we varied the fine-tuning rates from [0.01, 0.02, 0.05, 0.1], where the higher fine-tuning ratio, the more data that the defender can collect to conduct backdoor purification. Then, we vary the poisoning rate to simulate different adversary capability from [0.01, 0.02, 0.05, 0.1]. A defense should be stable and effective across the varied settings.

Effect of fine-tuning ratio. In this experiment, a larger fine-tuning ratio means a larger amount of data that the defender owns, while a small ratio is considered a more challenging setting. Figure 4 reports ASR (%) for two adaptive backdoors, COMBAT and SBL. The figure shows that the other baselines (FT, NAD, FST, TSBD) are highly sensitive to the fine-tuning ratio: their ASR reduction diminishes as the fine-tuning ratio decreases. Under COMBAT, these defenses still exhibit high ASR

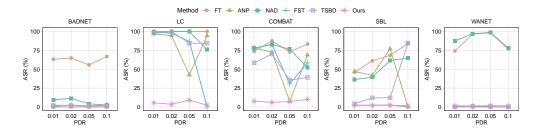


Figure 5: Defense results (ASR) under various poisoned data rate (PDR) settings with LC, COMBAT, and SBL attacks. The experiments are conducted on the CIFAR-10 dataset.

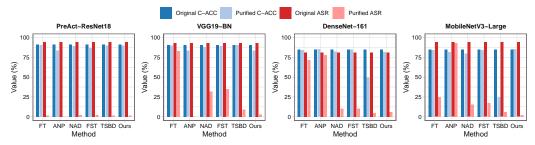


Figure 6: Defense results (ASR) under various model architecture settings with BadNet attack. The experiments are conducted on the CIFAR-10 dataset.

even at larger fine-tuning ratios. ANP can suppress ASR at favorable ratios but is unstable at smaller budgets. In contrast, UniBP achieves the lowest ASR across all ratios for both attacks, with the largest gains when the defender can use more than 2% of data for fine-tuning, highlighting superior sample efficiency and stability.

Effect of data poisoning rate. Figure 5 shows the effectiveness of all defenses versus different poisoned data rates (PDR) on CIFAR-10 across five attacks (BadNet, LC, Wanet, COMBAT, SBL). From the results, BadNet is the least challenging: all defenses achieve very low ASR. With LC, baselines are more sensitive to different poisoned data rates; several can only reduce ASR to 70% until when PDR is 0.05; whereas UnibP reduces ASR to as low as 0 across all ratios. Under COMBAT and SBL, even more advanced defenses such as TSBD and FST fluctuate widely and often exceed 50% even with larger budgets, and ANP is effective only at selective ratios. In contrast, UnibP is the most effective and stable across all attacks and fine-tuning ratios.

Analysis on model architecture. Figure 7 presents the effectiveness of different methods under various backbones: PreAct-ResNet18 (He et al., 2016), VGG19-BN (Simonyan & Zisserman, 2014), DenseNet-161 (Huang et al., 2017), and MobileNetV3-Large (Howard et al., 2019). We report both pre-defense (Original) and post-defense (Purified) clean accuracy (C-ACC) and attack success rate (ASR). Baseline fine-tuning defenses (FT, NAD, FST, TSBD) exhibit pronounced backbone dependence: on VGG19-BN and MobileNetV3-Large, they often leave high purified ASR or incur nontrivial C-ACC drops. ANP can substantially reduce ASR on some backbones (e.g., VGG19-BN) but typically at the cost of noticeable accuracy degradation. In contrast, UniBP consistently achieves the lowest ASR across all four architectures while keeping purified C-ACC close to the original, indicating model-agnostic effectiveness and a better robustness—accuracy trade-off.

5 Conclusion

We presented UniBP, a universal post-training defense for purifying backdoored models. The approach leverages BatchNorm statistics to expose backdoor footprints, rectifies these statistics on a small clean set, scores BN-affine parameters via a Fisher-based importance measure, prunes the most backdoor-sensitive entries, and fine-tunes with masked gradients—removing trigger pathways without prior knowledge of attack type or location. UniBP consistently attains the lowest ASR while preserving clean accuracy. It is stable across poisoning rates and fine-tuning budgets and operates effectively over a broad mask-ratio range, yielding strong robustness—accuracy trade-offs with modest clean data.

REFERENCES

- M. Barni, K. Kallas, and B. Tondi. A new backdoor attack in cnns by training set corruption without label poisoning. In *2019 IEEE International Conference on Image Processing (ICIP)*, pp. 101–105, 2019. doi: 10.1109/ICIP.2019.8802997.
- Bryant Chen, Wilka Carvalho, Nathalie Baracaldo, Heiko Ludwig, Benjamin Edwards, Taesung Lee, Ian Molloy, and Biplav Srivastava. Detecting backdoor attacks on deep neural networks by activation clustering, 2018. URL https://arxiv.org/abs/1811.03728.
- Weixin Chen, Baoyuan Wu, and Haoqian Wang. Effective backdoor defense by exploiting sensitivity of poisoned samples. *Advances in Neural Information Processing Systems*, 35:9727–9737, 2022.
- Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017.
- Pengzhou Cheng, Zongru Wu, Wei Du, Haodong Zhao, Wei Lu, and Gongshen Liu. Backdoor attacks and countermeasures in natural language processing models: A comprehensive security review. *IEEE Transactions on Neural Networks and Learning Systems*, 2025.
- Yansong Gao, Change Xu, Derui Wang, Shiping Chen, Damith C Ranasinghe, and Surya Nepal. Strip: A defence against trojan attacks on deep neural networks. In *Proceedings of the 35th annual computer security applications conference*, pp. 113–125, 2019.
- Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *ArXiv*, abs/1708.06733, 2017. URL https://api.semanticscholar.org/CorpusID:26783139.
- Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Evaluating backdooring attacks on deep neural networks. *IEEE Access*, 7:47230–47244, 2019.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European conference on computer vision*, pp. 630–645. Springer, 2016.
- Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, Quoc V. Le, and Hartwig Adam. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 1314–1324, 2019.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2261–2269, 2017.
- Kunzhe Huang, Yiming Li, Baoyuan Wu, Zhan Qin, and Kui Ren. Backdoor defense via decoupling the training process, 2022. URL https://arxiv.org/abs/2202.03423.
- Tran Huynh, Dang Nguyen, Tung Pham, and Anh Tran. Combat: Alternated training for effective clean-label backdoor attacks. In *AAAI Conference on Artificial Intelligence*, 2024. URL https://api.semanticscholar.org/CorpusID:268678332.
- Sagar Imambi, Kolla Bhanu Prakash, and GR Kanagachidambaresan. Pytorch. *Programming with TensorFlow: Solution for Edge Computing Applications*, pp. 87–104, 2021.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Yige Li, Nodens Koren, L. Lyu, Xixiang Lyu, Bo Li, and Xingjun Ma. Neural attention distillation: Erasing backdoor triggers from deep neural networks. *ArXiv*, abs/2101.05930, 2021a. URL https://api.semanticscholar.org/CorpusID:231627799.
 - Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Anti-backdoor learning: Training clean models on poisoned data. *Advances in Neural Information Processing Systems*, 34: 14900–14912, 2021b.

Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Neural attention distillation: Erasing backdoor triggers from deep neural networks. In *International Conference on Learning Representations*, 2021c. URL https://openreview.net/forum?id=910K4OM-OXE.

- Yige Li, Xixiang Lyu, Xingjun Ma, Nodens Koren, Lingjuan Lyu, Bo Li, and Yu-Gang Jiang. Reconstructive neuron pruning for backdoor defense. In *International Conference on Machine Learning*, pp. 19837–19854. PMLR, 2023a.
- Yiming Li, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. Backdoor learning: A survey. *IEEE transactions on neural networks and learning systems*, 35(1):5–22, 2022.
- Yudong Li, Shigeng Zhang, Weiping Wang, and Hong Song. Backdoor attacks to deep learning models and countermeasures: A survey. *IEEE Open Journal of the Computer Society*, 4:134–146, 2023b.
- Weilin Lin, Li Liu, Shaokui Wei, Jianze Li, and Hui Xiong. Unveiling and mitigating backdoor vulnerabilities based on unlearning weight changes and backdoor activeness. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=MfGRUVFtn9.
- Ibomoiye Domor Mienye and Theo G Swart. A comprehensive review of deep learning: Architectures, recent advances, and applications. *Information*, 15(12):755, 2024.
- Rui Min, Zeyu Qin, Li Shen, and Minhao Cheng. Towards stable backdoor purification through feature shift tuning. *ArXiv*, abs/2310.01875, 2023. URL https://api.semanticscholar.org/CorpusID:263608763.
- Rui Min, Zeyu Qin, Li Shen, and Minhao Cheng. Towards stable backdoor purification through feature shift tuning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Dung Thuy Nguyen, Ngoc N Tran, Taylor T Johnson, and Kevin Leach. Pbp: Post-training backdoor purification for malware classifiers. *arXiv preprint arXiv:2412.03441*, 2024.
- Tuan Anh Nguyen and Anh Tuan Tran. Wanet imperceptible warping-based backdoor attack. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=eEn8KTtJOx.
- Rui Ning, Jiang Li, Chunsheng Xin, and Hongyi Wu. Invisible poison: A blackbox clean label backdoor attack to deep neural networks. In *IEEE INFOCOM 2021 IEEE Conference on Computer Communications*, pp. 1–10, 2021. doi: 10.1109/INFOCOM42981.2021.9488902.
- Mohd Halim Mohd Noor and Ayokunle Olalekan Ige. A survey on state-of-the-art deep learning applications and challenges. *Engineering Applications of Artificial Intelligence*, 159:111225, 2025.
- Hoang Pham, The-Anh Ta, Anh Tran, and Khoa D. Doan. Flatness-aware sequential learning generates resilient backdoors. *ArXiv*, abs/2407.14738, 2024a. URL https://api.semanticscholar.org/CorpusID:271328781.
- Hoang Pham, The-Anh Ta, Anh Tran, and Khoa D Doan. Flatness-aware sequential learning generates resilient backdoors. In *European Conference on Computer Vision*, pp. 89–107. Springer, 2024b.
- Wojciech Samek, Grégoire Montavon, Sebastian Lapuschkin, Christopher J Anders, and Klaus-Robert Müller. Explaining deep neural networks and beyond: A review of methods and applications. *Proceedings of the IEEE*, 109(3):247–278, 2021.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Johannes Stallkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. The German Traffic Sign Recognition Benchmark: A multi-class classification competition. In *IEEE International Joint Conference on Neural Networks*, pp. 1453–1460, 2011.

- Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks. *arXiv* preprint arXiv:1912.02771, 2019.
 - Yichen Wan, Youyang Qu, Wei Ni, Yong Xiang, Longxiang Gao, and Ekram Hossain. Data and model poisoning backdoor attacks on wireless federated learning, and the defense mechanisms: A comprehensive survey. *IEEE Communications Surveys & Tutorials*, 26(3):1861–1897, 2024.
 - Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In 2019 IEEE symposium on security and privacy (SP), pp. 707–723. IEEE, 2019.
 - Hang Wang, Zhen Xiang, David J. Miller, and George Kesidis. Mm-bd: Post-training detection of backdoor attacks with arbitrary backdoor pattern types using a maximum margin statistic, 2023. URL https://arxiv.org/abs/2205.06900.
 - Zhenting Wang, Juan Zhai, and Shiqing Ma. Bppattack: Stealthy and efficient trojan attacks against deep neural networks via image quantization and contrastive adversarial learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15074–15084, June 2022.
 - Baoyuan Wu, Hongrui Chen, Mingda Zhang, Zihao Zhu, Shaokui Wei, Danni Yuan, and Chao Shen. Backdoorbench: A comprehensive benchmark of backdoor learning. *Advances in Neural Information Processing Systems*, 35:10546–10559, 2022.
 - Dongxian Wu and Yisen Wang. Adversarial neuron pruning purifies backdoored deep models. *Advances in Neural Information Processing Systems*, 34:16913–16925, 2021.
 - Yi Zeng, Si Chen, Won Park, Zhuoqing Mao, Ming Jin, and Ruoxi Jia. Adversarial unlearning of backdoors via implicit hypergradient. In *International Conference on Learning Representations*, 2021.
 - Yi Zeng, Si Chen, Won Park, Zhuoqing Mao, Ming Jin, and Ruoxi Jia. Adversarial unlearning of backdoors via implicit hypergradient. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=MeeQkFYVbzW.
 - Yi Zeng, Minzhou Pan, Hoang Anh Just, Lingjuan Lyu, Meikang Qiu, and Ruoxi Jia. Narcissus: A practical clean-label backdoor attack with limited information. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, pp. 771–785, 2023.
 - Shaobo Zhang, Yimeng Pan, Qin Liu, Zheng Yan, Kim-Kwang Raymond Choo, and Guojun Wang. Backdoor attacks and defenses targeting multi-domain ai models: A comprehensive review. *ACM Computing Surveys*, 57(4):1–35, 2024.
 - Runkai Zheng, Rongjun Tang, Jianze Li, and Li Liu. Data-free backdoor removal based on channel lipschitzness. In *European Conference on Computer Vision*, pp. 175–191. Springer, 2022a.
 - Runkai Zheng, Rongjun Tang, Jianze Li, and Li Liu. Pre-activation distributions expose backdoor neurons. *Advances in Neural Information Processing Systems*, 35:18667–18680, 2022b.
 - Mingli Zhu, Shaokui Wei, Li Shen, Yanbo Fan, and Baoyuan Wu. Enhancing fine-tuning based backdoor defense with sharpness-aware minimization. 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 4443–4454, 2023. URL https://api.semanticscholar.org/CorpusID: 258297949.

APPENDIX

We conduct all the experiments using PyTorch 2.1.0 (Imambi et al., 2021). All experiments are run on a computer with an Intel Xeon Gold 6330N CPU and an NVIDIA A6000 GPU.

A DETAILED EXPERIMENTAL SETUP

A.1 DATASETS AND PREPROCESSING

CIFAR-10. The CIFAR-10 dataset (Krizhevsky et al., 2009) comprises $60,000\ 32\times32$ RGB images evenly distributed across 10 classes. We adopt the official split with 50,000 training images and 10,000 test images (6,000 per class in total; 5,000 train and 1,000 test per class). Unless otherwise noted, we follow the standard evaluation protocol on the test set.

GTSRB (German Traffic Sign Recognition Benchmark). The GTSRB dataset (Stallkamp et al., 2011) contains 51,839 images across 43 classes, with 39,209 images for training and 12,630 for testing. Following common practice, we use the standard train/test split and resize all images to 32×32 RGB for training and evaluation.

A.2 ATTACK DETAILS

We evaluate our defense against five SOTA backdoor attacks: BadNets (Gu et al., 2019), LC (Turner et al., 2019), WaNet (Nguyen & Tran, 2021), COMBAT (Huynh et al., 2024), and SBL (Pham et al., 2024b). For BadNets, LC, and WaNet, we adopt the implementations provided in the BackdoorBench framework and use the default configurations. Since COMBAT and SBL are not integrated into BackdoorBench, we incorporated them into our codebase using the official implementations released by the authors¹ 2. To ensure consistency and comparability across all experiments, we fixed the poisoning ratio at 10%. Examples of poisoned images under each attack is shown in Figure 8.

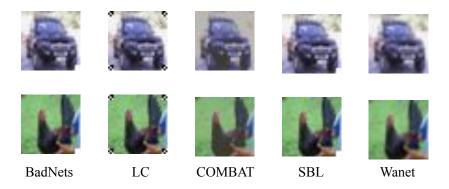


Figure 7: Examples of poisoned images on the CIFAR-10 dataset.

A.3 DEFENSE DETAILS

We compare our method against five SOTA backdoor defenses: FT, ANP (Wu & Wang, 2021), NAD (Li et al., 2021a), FST (Min et al., 2024), and TSBD (Lin et al., 2024). For FT, ANP, NAD, and TSBD, we adopt the implementations and default configurations provided in the BackdoorBench framework. Since FST is not included in BackdoorBench, we integrated it into our codebase using the publicly released implementation ³. To ensure fairness across methods, we set the batch size to 256 for all defenses, except for FST, where we follow the original paper and use a batch size of 128.

https://github.com/VinAIResearch/COMBAT

https://github.com/mail-research/SBL-resilient-backdoors

https://github.com/AISafety-HKUST/stable_backdoor_purification

All defenses are trained with a learning rate of 0.002 for 20 epochs. For TSBD, we follow the settings reported in the original paper, fixing the neuron ratio at n = 0.15 and the weight ratio at m = 0.7.

A.4 MODEL ARCHITECTURES AND INITIALIZATION

We evaluate four backbone architectures representative of common vision families:

PreAct-ResNet-18. Standard PreAct-ResNet-18; the final classifier is replaced to match the dataset classes (10 for CIFAR-10; 43 for GTSRB).

VGG19-BN. VGG19 with batch normalization after each convolutional block; initialized from ImageNet and refit with a dataset-specific classifier.

DenseNet-161. ImageNet-pretrained DenseNet-161; the classifier head is replaced to match the target classes.

MobileNetV3-Large. ImageNet-pretrained MobileNetV3-Large; the final fully connected layer is replaced to fit the dataset classes.

Model modifications for purification: Our purification pipeline interacts primarily with BatchNorm affine parameters and per-channel statistics. We instrument BatchNorm layers to read and optionally reset γ, β and moving averages ($\mu_{\rm mov}, \sigma_{\rm mov}$). For the pruning / affine-mask step we add small, lightweight selection masks per channel (implemented as binary or continuous gates) that can be applied to the BN affine scale term γ during inference and finetuning.

A.5 HYPER-PARAMETERS

The pipeline includes separate hyper-parameters for (A) initial training/victim model creation (poisoned model), and (B) purification stages. We list the values used in all experiments unless noted otherwise.

Training Phase. Unless otherwise noted, poisoned models are trained using PreAct-ResNet-18 with SGD (momentum 0.9), an initial learning rate of 0.01, weight decay of 5×10^{-4} , batch size 128, and 100 epochs. The learning rate follows *CosineAnnealingLR*. The random seed is fixed to 0. Unless otherwise specified, standard data augmentation (random horizontal flip and random crop) is applied.

Fine-tuning Phase. During fine-tuning, we use SGD with momentum 0.9 and a learning rate in the range 1×10^{-3} to 2×10^{-4} ; unless otherwise specified, the training batch size is 128. In the sensitivity-to-fine-tuning-ratio study, we sweep the fine-tuning ratio over $\{1\%, 2\%, 5\%, 10\%\}$ and adjust batch sizes accordingly, i.e., training mini-batch is $\{32, 32, 64, 128\}$, respectively.

B ADDITIONAL EXPERIMENTAL RESULTS

B.1 SENSITIVITY TO DATA POISONING RATES

Table 3 presents a comprehensive evaluation of different defense methods (Pretrained, FT, ANP, NAD, FST, TSBD, and UnibP) against a range of backdoor attacks including BadNet, LC, COMBAT, SBL, and Wanet, under varying poisoning data ratios (PDR = 0.1, 0.05, 0.02, 0.01). For each configuration, both model accuracy (MA) and attack success rate (ASR) are reported to highlight the trade-off between maintaining clean accuracy and suppressing malicious behavior. Across the board, baseline Pretrained models show high MA but consistently elevated ASR, indicating vulnerability to all attacks. Fine-tuning (FT) improves resilience to some extent, though it struggles to reduce ASR under low PDRs. ANP and NAD demonstrate stronger backdoor mitigation, often reducing ASR close to zero, but at the cost of a slight drop in MA in some cases. FST and TSBD provide a more balanced trade-off, achieving high MA while substantially lowering ASR in multiple attack settings. Notably, UniBP consistently achieves competitive MA while keeping ASR at very low levels, especially under LC and COMBAT attacks, showcasing its robustness under challenging conditions. Overall, the results emphasize that while most defenses reduce ASR to some degree, methods like NAD, FST, TSBD, and particularly UniBP stand out in delivering both strong protection and reliable utility.

Table 3: Performance comparison of different defense methods (Pretrained, FT, ANP, NAD, FST, TSBD, and PBP) against multiple backdoor attacks (BadNet, LC, COMBAT, SBL, and Wanet) under varying poisoning data ratios (PDR = 0.1, 0.05, 0.02, 0.01). The table reports the model accuracy (MA) and attack success rate (ASR) in percentage.

| Attacks | | Pretrained | | FT | | ANP | | NAD | | FST | | TSBD | | Ours | |
|---------|------|------------|--------|-------|--------|-------|-------|-------|--------|-------|-------|-------|-------|-------|-------|
| | PDR | C-ACC | ASR | C-ACC | ASR | C-ACC | ASR | C-ACC | ASR | C-ACC | ASR | C-ACC | ASR | C-ACC | ASR |
| BadNet | 0.1 | 91.44 | 94.41 | 91.01 | 66.94 | 84.10 | 0.00 | 89.62 | 2.66 | 91.48 | 3.13 | 91.67 | 2.10 | 89.82 | 1.47 |
| | 0.05 | 92.15 | 90.30 | 91.37 | 55.80 | 85.05 | 0.00 | 90.58 | 4.24 | 91.37 | 1.30 | 92.19 | 1.43 | 87.85 | 1.50 |
| | 0.02 | 92.81 | 81.47 | 91.77 | 64.78 | 86.49 | 0.01 | 91.34 | 11.52 | 92.33 | 2.58 | 92.18 | 1.88 | 87.59 | 2.09 |
| | 0.01 | 93.34 | 71.21 | 92.39 | 63.20 | 85.05 | 0.01 | 91.12 | 9.63 | 92.50 | 0.94 | 93.02 | 1.94 | 87.68 | 2.51 |
| LC | 0.1 | 84.19 | 100.00 | 92.80 | 100.00 | 84.67 | 94.48 | 91.39 | 75.88 | 91.08 | 0.00 | 91.08 | 84.27 | 89.09 | 2.36 |
| | 0.05 | 93.32 | 100.00 | 92.26 | 100.00 | 91.55 | 42.47 | 90.82 | 100.00 | 92.22 | 86.91 | 91.66 | 84.52 | 86.19 | 9.13 |
| | 0.02 | 93.39 | 100.00 | 92.68 | 100.00 | 84.67 | 94.48 | 91.57 | 99.91 | 92.38 | 97.55 | 92.98 | 99.84 | 83.06 | 3.67 |
| | 0.01 | 93.54 | 99.97 | 92.36 | 100.00 | 88.77 | 96.80 | 91.51 | 99.01 | 92.59 | 99.52 | 92.78 | 97.54 | 83.04 | 5.61 |
| | 0.1 | 85.05 | 99.23 | 91.90 | 83.47 | 84.20 | 68.90 | 91.80 | 52.28 | 92.27 | 58.40 | 92.56 | 39.20 | 91.04 | 10.28 |
| COMBAT | 0.05 | 93.94 | 94.47 | 93.46 | 72.83 | 85.18 | 7.58 | 98.41 | 76.56 | 91.25 | 30.65 | 92.91 | 35.57 | 87.90 | 7.18 |
| COMBAI | 0.02 | 93.90 | 85.04 | 94.20 | 87.63 | 85.05 | 72.56 | 93.61 | 82.51 | 93.90 | 82.73 | 92.76 | 70.21 | 89.72 | 6.23 |
| | 0.01 | 94.14 | 83.67 | 93.49 | 73.76 | 93.40 | 78.12 | 93.40 | 78.12 | 93.60 | 78.07 | 92.78 | 58.46 | 87.29 | 7.72 |
| | 0.1 | 90.52 | 88.84 | 90.79 | 83.85 | 88.77 | 0.04 | 90.39 | 64.80 | 91.17 | 0.24 | 91.43 | 84.68 | 83.25 | 1.86 |
| SBL | 0.05 | 90.02 | 79.35 | 89.94 | 68.31 | 85.53 | 77.16 | 89.68 | 61.60 | 90.06 | 2.58 | 89.55 | 12.15 | 87.81 | 2.33 |
| SDL | 0.02 | 90.25 | 68.27 | 90.33 | 61.16 | 87.59 | 41.98 | 90.03 | 39.53 | 90.17 | 2.48 | 90.25 | 12.12 | 88.04 | 2.13 |
| | 0.01 | 90.50 | 47.07 | 90.39 | 46.14 | 90.50 | 47.07 | 89.98 | 36.30 | 90.49 | 2.08 | 91.01 | 4.33 | 88.44 | 2.60 |
| | 0.1 | 93.40 | 99.97 | 93.75 | 76.73 | 84.46 | 0.08 | 93.72 | 78.07 | 93.04 | 0.30 | 77.64 | 0.00 | 90.22 | 1.52 |
| Wanet | 0.05 | 93.37 | 99.87 | 93.89 | 98.14 | 85.33 | 0.01 | 93.88 | 98.73 | 93.32 | 0.27 | 70.42 | 0.70 | 90.09 | 1.88 |
| wanet | 0.02 | 93.43 | 99.38 | 93.54 | 96.54 | 88.44 | 0.13 | 93.77 | 96.98 | 93.37 | 0.50 | 43.97 | 1.15 | 89.47 | 1.39 |
| | 0.01 | 93.81 | 98.48 | 93.80 | 74.12 | 85.88 | 0.21 | 93.83 | 87.26 | 93.18 | 1.13 | 75.40 | 0.00 | 89.71 | 1.73 |

Table 4: Comparison of defenses against BadNet across models (original vs. purified). Values are clean accuracy (C-ACC) and attack success rate (ASR), both in %.

| Model | Tag | FT | | ANP | | NAD | | FST | | TSBD | | Ours | |
|-------------------|----------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1710401 | 145 | C-ACC | ASR |
| VGG19_BN | Original | 90.84 | 93.41 | 90.84 | 93.41 | 90.84 | 93.41 | 90.84 | 93.41 | 90.84 | 93.41 | 90.84 | 93.41 |
| | Purified | 89.72 | 82.77 | 83.38 | 0.00 | 88.09 | 32.15 | 89.05 | 34.74 | 90.67 | 9.38 | 83.48 | 2.84 |
| DenseNet161 | Original | 84.97 | 80.91 | 84.97 | 80.91 | 84.97 | 80.91 | 84.97 | 80.91 | 84.97 | 80.91 | 84.97 | 80.91 |
| | Purified | 84.58 | 71.77 | 85.28 | 78.18 | 82.51 | 10.30 | 84.87 | 10.13 | 49.72 | 4.88 | 82.03 | 5.70 |
| MobileNetV3-Large | Original | 85.12 | 94.12 | 85.12 | 94.12 | 85.12 | 94.12 | 85.12 | 94.12 | 85.12 | 94.12 | 85.12 | 94.12 |
| | Purified | 84.24 | 24.71 | 81.74 | 93.36 | 79.80 | 15.62 | 84.20 | 17.15 | 25.13 | 5.86 | 85.78 | 2.23 |
| PreAct-ResNet18 | Original | 91.44 | 94.41 | 91.44 | 94.41 | 91.44 | 94.41 | 91.44 | 94.41 | 91.44 | 94.41 | 91.44 | 94.41 |
| | Purified | 90.56 | 1.47 | 83.51 | 0.00 | 89.33 | 2.08 | 87.06 | 2.08 | 90.13 | 1.78 | 89.82 | 1.47 |

B.2 Sensitivity to model architectures

Table 4 compares Original vs. Purified C-ACC/ASR under BadNet across four backbones. Baseline fine-tuning defenses (FT, NAD, FST, TSBD) show pronounced backbone dependence: on VGG19-BN and DenseNet161 they often leave high purified ASR ($\sim 10\%$), and on DenseNet161 and MobileNetV3-Large, TSBD substantially reduces C-ACC. ANP lowers ASR on some backbones (VGG19-BN, PreAct-ResNet18) but with noticeable accuracy drops (7-8%) and fails on the others (ASR remains high, often above 70-90%). In contrast, UniBP keeps ASR low across all architectures (about 1-6%) while maintaining purified C-ACC close to the original (typically within a few points), indicating backbone-agnostic effectiveness and a better robustness—accuracy trade-off.

B.3 ABLATION STUDY

We sweep the mask ratio K, the primary control in our method, and summarize the outcomes in Figure 8. Across all settings, C-ACC decreases smoothly as K increases, with only a small drop (typically ≤ 5 points) inside the shaded range and a sharp decline once $K \geq 0.10 \times 10^{-3}$. ASR remains low overall, generally within 1–5%; LC at 10% poisoning shows a mild bump near $K \approx 0.06 \times 10^{-3}$, but the trend is otherwise flat. Increasing K beyond the shaded range yields little additional ASR reduction while causing substantial loss in clean accuracy, most notably for BadNet at 5% poisoning. Small pruning budgets within the highlighted range therefore, provide the best trade-off, keeping ASR low with minimal impact on clean performance across both attack families and poisoning rates.

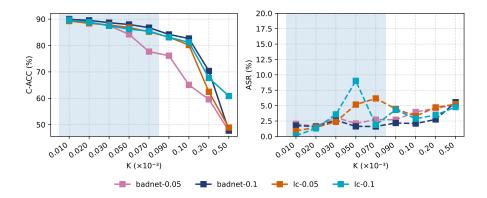


Figure 8: Effect of pruning budget K on clean accuracy (C-ACC, left) and attack success rate (ASR, right) under BadNet and LC with poisoning rates 5% and 10% on CIFAR-10. The shaded band marks the stable operating range ($K \in [0.010, 0.070] \times 10^{-3}$).

B.4 Additional Plots

 Figure 9 summarizes how different backdoor families distort the representation space and BatchNorm statistics. The t-SNE plots (top) show that BadNet and LC largely blend poisoned samples into the target-class manifold, yielding only mild geometric separation; WANET induces a moderate shift with partially segregated clusters; SBL creates a compact, outlying poisoned cluster that is clearly detached from clean structure; COMBAT, which mixes patch- and distributional cues, produces overlap similar to BadNet but with denser target-class concentration. The histograms of BN perchannel means (bottom) mirror these trends: BadNet and LC exhibit near-overlapping clean vs. backdoored distributions (small mean shifts), WANET shows a visible but modest shift, and SBL displays a pronounced displacement of the backdoored distribution. COMBAT lies between these extremes. Overall, attacks that strongly perturb intermediate distributions (e.g., SBL) leave a larger BN footprint, whereas patch-like attacks (BadNet/LC) are more stealthy in BN space—motivating a rectification objective that leverages BN statistics while also requiring parameter-level masking to handle the subtler cases. Though these attacks are different in manner and how the trigger is crafted, the shift phenomenon in BN statistics could be leveraged to defend against these attacks.

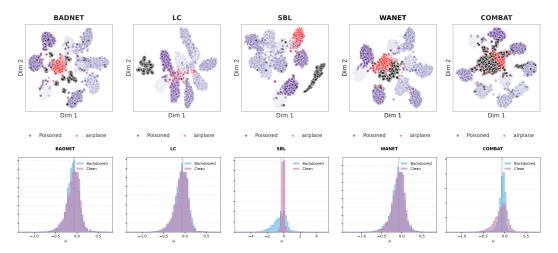


Figure 9: t-SNE of feature embeddings of different attack strategies and their effect on BN layers' statistic CIFAR-10 of different attack families.

C LIMITATIONS

 We note several limitations that contextualize our results and suggest directions for future work. First, the method assumes access to a small hold-out clean set to estimate BatchNorm statistics and to drive affine-mask learning; its size, class coverage, and label quality materially affect stability and final accuracy. In extremely low-data or noisy-label regimes, the rectification signal can weaken, and the fully unsupervised setting (no clean data) is outside our scope. Second, while we evaluate adaptive variants, a stronger adversary that co-designs triggers to survive BN-affine reset and pruning, perturbs or hijacks running statistics during poisoning, or disperses triggers to reduce gradient salience could diminish effectiveness; developing defenses with explicit guarantees against such adaptive strategies remains open. Third, our study focuses on image classification with BN-based architectures; extending the approach to other modalities (e.g., audio, NLP) or tasks (e.g., detection, segmentation), and to models using alternative normalizations (e.g., LayerNorm, GroupNorm), will require adapting both the rectification objective and the mask parameterization.

D BROADER IMPACT

Positive impacts. The method strengthens deployed classifiers against poisoning/backdoor threats, improving robustness in safety-critical settings (e.g., automotive perception, medical imaging).

Dual use. Defensive techniques can inform stronger, defense-aware attacks. We will release code with clear usage guidance and a responsible license, and provide deployment recommendations (e.g., separate clean validation, periodic re-evaluation), limiting exploit-ready details to what is necessary for reproducibility.

Privacy. The approach assumes a small clean dataset; when data are sensitive, practitioners should minimize collection, de-identify inputs, restrict access, and follow IRB requirements.

Responsible disclosure. We support coordinated disclosure to affected stakeholders and commit to sharing only information needed for verification and remediation.